

Perceptron and Logistic Regression

Jindřich Libovický (reusing materials by Milan Straka)

 October 17, 2023



Charles University in Prague
Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics



unless otherwise stated

After this lecture you should be able to

- Think about binary classification using **geometric intuition** and use the **perceptron algorithm**.
- Define the **main concepts of information theory** (entropy, cross-entropy, KL-divergence) and prove their properties.
- Derive training objectives using the **maximum likelihood principle**.
- Implement and use **logistic regression** for binary classification with SGD.

Binary Classification

Binary classification is a classification in two classes.

The simplest way to evaluate classification is **accuracy**, which is the ratio of input examples that were classified correctly – i.e., where the predicted class and the target class match.

To extend linear regression to binary classification, we might seek a **threshold** and then classify an input as negative/positive depending on whether $y(\mathbf{x}; \mathbf{w}) = \mathbf{x}^T \mathbf{w} + b$ is smaller/larger than a given threshold.

Zero value is usually used as the threshold, both because of symmetry and also because the **bias** parameter acts as a trainable threshold anyway.

The set of points with prediction 0 is called a **decision boundary**.

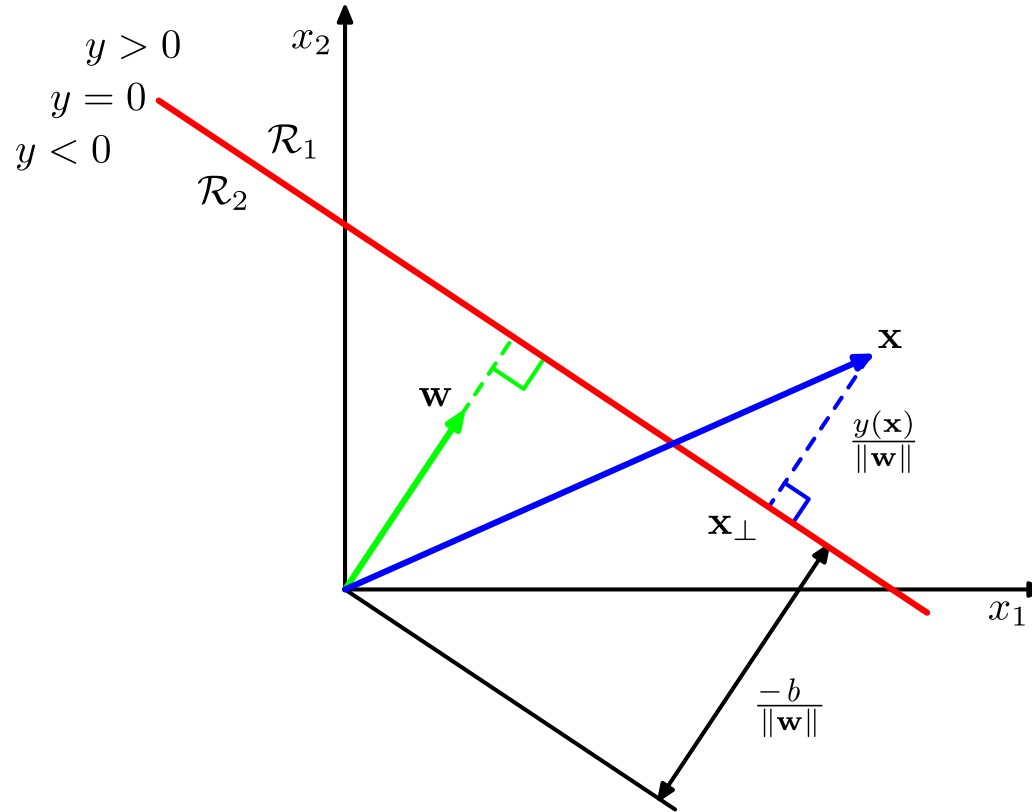


Figure 4.1 of Pattern Recognition and Machine Learning.

The perceptron algorithm is probably the oldest one for training weights of a binary classification. Assuming the target value $t \in \{-1, +1\}$, the goal is to find weights \mathbf{w} such that for all train data,

$$\text{sign}(y(\mathbf{x}_i; \mathbf{w})) = \text{sign}(\mathbf{x}_i^T \mathbf{w}) = t_i,$$

or equivalently,

$$t_i y(\mathbf{x}_i; \mathbf{w}) = t_i \mathbf{x}_i^T \mathbf{w} > 0.$$

Note that a set is called **linearly separable**, if there exists a weight vector \mathbf{w} such that the above equation holds.

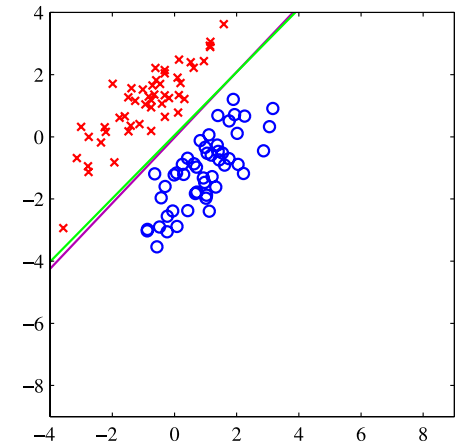


Figure 4.4 of Pattern Recognition and Machine Learning.

The perceptron algorithm was invented by Rosenblatt in 1958.

Input: Linearly separable dataset ($\mathbf{X} \in \mathbb{R}^{N \times D}$, $\mathbf{t} \in \{-1, +1\}^N$).

Output: Weights $\mathbf{w} \in \mathbb{R}^D$ such that $t_i \mathbf{x}_i^T \mathbf{w} > 0$ for all i .

- $\mathbf{w} \leftarrow \mathbf{0}$
- until all examples are classified correctly, process example i :
 - $y \leftarrow \mathbf{x}_i^T \mathbf{w}$
 - if $t_i y \leq 0$ (incorrectly classified example):
 - $\mathbf{w} \leftarrow \mathbf{w} + t_i \mathbf{x}_i$

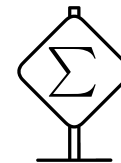
We will prove that the algorithm always arrives at some correct set of weights \mathbf{w} if the training set is linearly separable.

Proof of Perceptron Convergence

Let \mathbf{w}_* be some weights correctly classifying (separating) the training data, and let \mathbf{w}_k be the weights after k nontrivial updates of the perceptron algorithm, with \mathbf{w}_0 being 0.

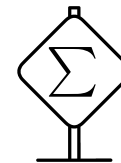
We will prove that the angle α between \mathbf{w}_* and \mathbf{w}_k decreases at each step. Note that

$$\cos(\alpha) = \frac{\mathbf{w}_*^T \mathbf{w}_k}{\|\mathbf{w}_*\| \cdot \|\mathbf{w}_k\|}.$$



Proof of Perceptron Convergence

Assume that the maximum norm of any training example $\|\mathbf{x}\|$ is bounded by R , and that γ is the minimum margin of \mathbf{w}_* , so for each training example (\mathbf{x}, t) , $t\mathbf{x}^T \mathbf{w}_* \geq \gamma$.



First consider the dot product of \mathbf{w}_* and \mathbf{w}_k :

$$\mathbf{w}_*^T \mathbf{w}_k = \mathbf{w}_*^T (\mathbf{w}_{k-1} + t_k \mathbf{x}_k) \geq \mathbf{w}_*^T \mathbf{w}_{k-1} + \gamma.$$

By iteratively applying this equation, we get

$$\mathbf{w}_*^T \mathbf{w}_k \geq k\gamma.$$

Now consider the length of \mathbf{w}_k :

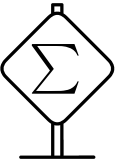
$$\|\mathbf{w}_k\|^2 = \|\mathbf{w}_{k-1} + t_k \mathbf{x}_k\|^2 = \|\mathbf{w}_{k-1}\|^2 + 2t_k \mathbf{x}_k^T \mathbf{w}_{k-1} + \|\mathbf{x}_k\|^2.$$

Because \mathbf{x}_k was misclassified, we know that $t_k \mathbf{x}_k^T \mathbf{w}_{k-1} \leq 0$, so $\|\mathbf{w}_k\|^2 \leq \|\mathbf{w}_{k-1}\|^2 + R^2$.

When applied iteratively, we get $\|\mathbf{w}_k\|^2 \leq k \cdot R^2$.

Putting everything together, we get

$$\cos(\alpha) = \frac{\mathbf{w}_*^T \mathbf{w}_k}{\|\mathbf{w}_*\| \cdot \|\mathbf{w}_k\|} \geq \frac{k\gamma}{\sqrt{k}R^2\|\mathbf{w}_*\|}.$$



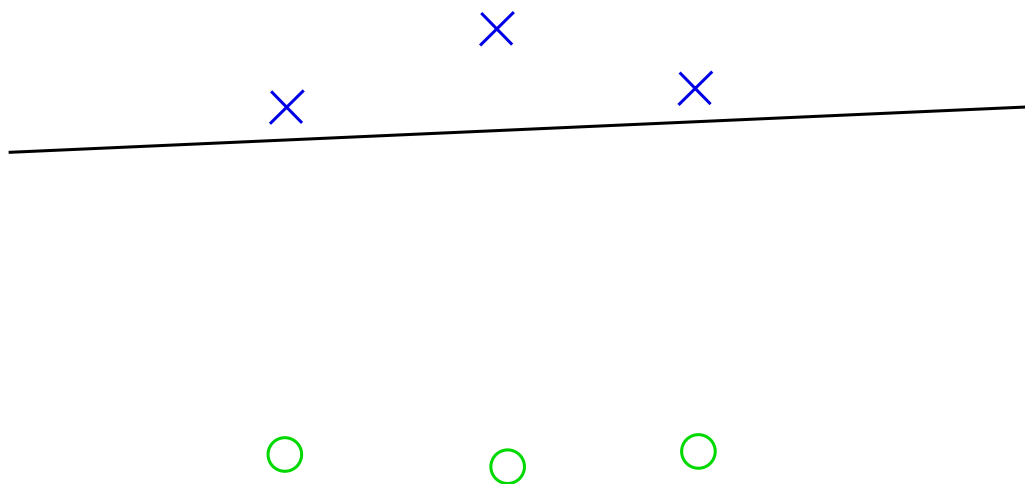
Therefore, the $\cos(\alpha)$ increases during every update. Because the value of $\cos(\alpha)$ is at most one, we can compute the upper bound on the number of steps when the algorithm converges as

$$1 \geq \frac{\sqrt{k}\gamma}{\sqrt{R^2}\|\mathbf{w}_*\|} \text{ or } k \leq \frac{R^2\|\mathbf{w}_*\|^2}{\gamma^2}.$$

Perceptron Issues

Perceptron has several drawbacks:

- If the input set is not linearly separable, the algorithm never finishes.
- The algorithm performs only prediction, it is not able to return the probabilities of predictions.
- Most importantly, Perceptron algorithm finds *some* solution, not necessarily a good one, because once it finds some, it cannot perform any more updates.



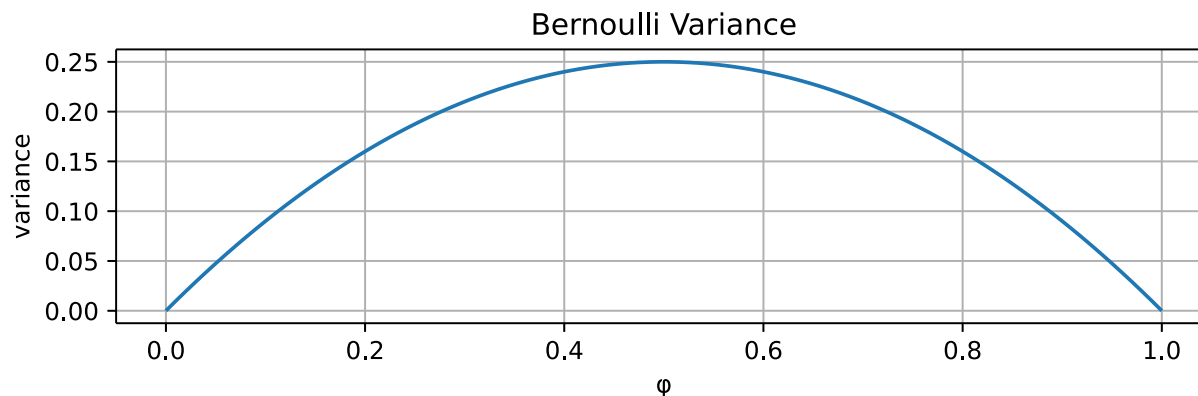
Bernoulli Distribution

The Bernoulli distribution is a distribution over a binary random variable. It has a single parameter $\varphi \in [0, 1]$, which specifies the probability that the random variable is equal to 1.

$$P(x) = \varphi^x (1 - \varphi)^{1-x}$$

$$\mathbb{E}[x] = \varphi$$

$$\text{Var}(x) = \varphi(1 - \varphi)$$



Categorical Distribution

Extension of the Bernoulli distribution to random variables taking one of K different discrete outcomes. It is parametrized by $\mathbf{p} \in [0, 1]^K$ such that $\sum_{i=0}^{K-1} p_i = 1$.

We represent outcomes as vectors $\in \{0, 1\}^K$ in the **one-hot encoding**. Therefore, an outcome $x \in \{0, 1, \dots, K - 1\}$ is represented as a vector

$$\mathbf{1}_x \stackrel{\text{def}}{=} \left([i = x] \right)_{i=0}^{K-1} = \left(\underbrace{0, \dots, 0}_x, 1, \underbrace{0, \dots, 0}_{K-x-1} \right).$$

The outcome probability, mean, and variance are very similar to the Bernoulli distribution.

$$\begin{aligned} P(\mathbf{x}) &= \prod_{i=0}^{K-1} p_i^{x_i} \\ \mathbb{E}[x_i] &= p_i \\ \text{Var}(x_i) &= p_i(1 - p_i) \end{aligned}$$

Self Information

Amount of **surprise** when a random variable is sampled.

- Should be zero for events with probability 1.
- Less likely events are more surprising.
- Independent events should have **additive** information.

$$I(x) \stackrel{\text{def}}{=} -\log P(x) = \log \frac{1}{P(x)}$$

Entropy

Amount of **surprise** in the whole distribution.

$$H(P) \stackrel{\text{def}}{=} \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\log P(x)]$$

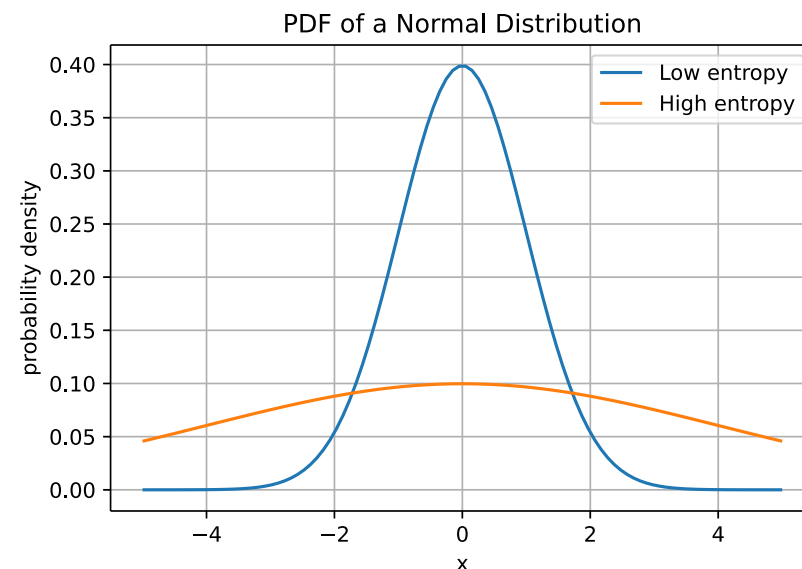
- for discrete P : $H(P) = -\sum_x P(x) \log P(x)$
- for continuous P : $H(P) = -\int P(x) \log P(x) dx$

Because $\lim_{x \rightarrow 0} x \log x = 0$, for $P(x) = 0$ we consider $P(x) \log P(x)$ to be zero.

Note that in the continuous case, the continuous entropy (also called *differential entropy*) has slightly different semantics, for example, it can be negative.

For binary logarithms, the entropy is measured in **bits**.

However, from now on, all logarithms are *natural logarithms* with base e (and then the entropy is measured in units called **nats**).



Cross-Entropy

$$H(P, Q) \stackrel{\text{def}}{=} -\mathbb{E}_{x \sim P} [\log Q(x)]$$

Gibbs Inequality

- $H(P, Q) \geq H(P)$
- $H(P) = H(P, Q) \Leftrightarrow P = Q$

Proof: Consider $H(P) - H(P, Q) = \sum_x P(x) \log \frac{Q(x)}{P(x)}$.

Using the fact that $\log x \leq (x - 1)$ with equality only for $x = 1$, we get

$$\sum_x P(x) \log \frac{Q(x)}{P(x)} \leq \sum_x P(x) \left(\frac{Q(x)}{P(x)} - 1 \right) = \sum_x Q(x) - \sum_x P(x) = 0.$$

For the equality to hold, $\frac{Q(x)}{P(x)}$ must be 1 for all x , i.e., $P = Q$.

Kullback-Leibler Divergence (KL Divergence)

Sometimes also called **relative entropy**.

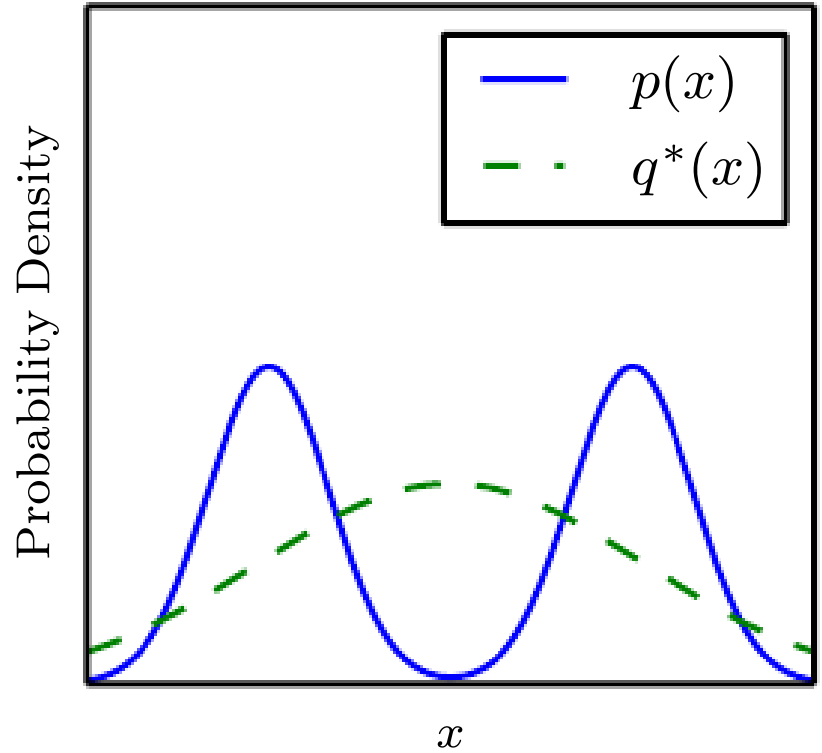
$$D_{\text{KL}}(P\|Q) \stackrel{\text{def}}{=} H(P, Q) - H(P) = \mathbb{E}_{x \sim P}[\log P(x) - \log Q(x)]$$

- consequence of Gibbs inequality: $D_{\text{KL}}(P\|Q) \geq 0$, $D_{\text{KL}}(P\|Q) = 0$ iff $P = Q$
- generally $D_{\text{KL}}(P\|Q) \neq D_{\text{KL}}(Q\|P)$

Nonsymmetry of KL Divergence

Assume we want find the best unimodal distribution q .

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p||q)$$



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q||p)$$

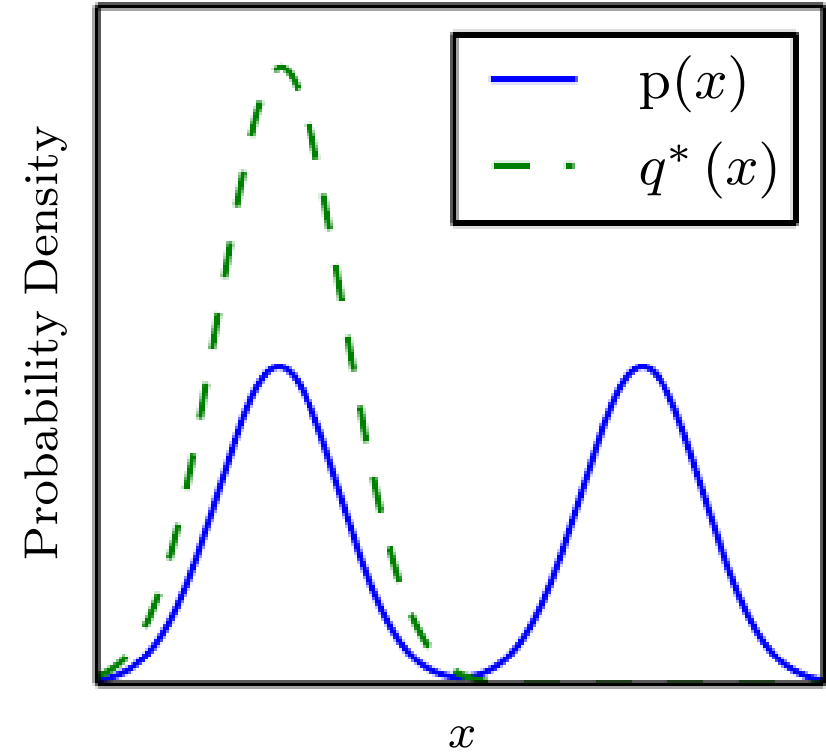


Figure 3.6 of "Deep Learning" book, <https://www.deeplearningbook.org>

Normal (or Gaussian) Distribution

Distribution over real numbers, parametrized by a mean μ and variance σ^2 :

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

For standard values $\mu = 0$ and $\sigma^2 = 1$ we get $\mathcal{N}(x; 0, 1) = \sqrt{\frac{1}{2\pi}} e^{-\frac{x^2}{2}}$.

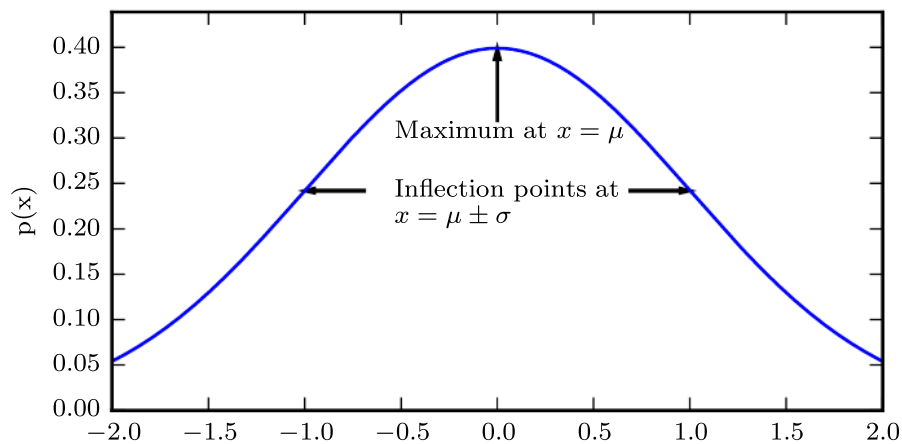


Figure 3.1 of "Deep Learning" book, <https://www.deeplearningbook.org>.

Central Limit Theorem

The sum of independent identically distributed random variables with finite variance converges to normal distribution.

Principle of Maximum Entropy

Given a set of constraints, a distribution with maximal entropy fulfilling the constraints can be considered the most general one, containing as little additional assumptions as possible.

Considering distributions with a **given mean and variance**, it can be proven (using variational inference) that such a distribution with **maximum entropy** is exactly the normal distribution.

Maximum Likelihood Estimation

Let $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ be training data drawn independently from the data-generating distribution p_{data} .

We denote the **empirical data distribution** as \hat{p}_{data} , where

$$\hat{p}_{\text{data}}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{|\{i : \mathbf{x}_i = \mathbf{x}\}|}{N}.$$

Let $p_{\text{model}}(\mathbf{x}; \mathbf{w})$ be a family of distributions.

- If the weights are fixed, $p_{\text{model}}(\mathbf{x}; \mathbf{w})$ is a probability distribution.
- If we instead consider the fixed training data \mathbf{X} , then

$$L(\mathbf{w}) = p_{\text{model}}(\mathbf{X}; \mathbf{w}) = \prod_{i=1}^N p_{\text{model}}(\mathbf{x}_i; \mathbf{w})$$

is called the **likelihood**. Note that even if the value of the likelihood is in range $[0, 1]$, it is not a probability, because the likelihood is not a probability distribution.

Maximum Likelihood Estimation

Let $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ be training data drawn independently from the data-generating distribution p_{data} . We denote the empirical data distribution as \hat{p}_{data} and let $p_{\text{model}}(\mathbf{x}; \mathbf{w})$ be a family of distributions.

The **maximum likelihood estimation** of \mathbf{w} is:

$$\begin{aligned}
 \mathbf{w}_{\text{MLE}} &= \arg \max_{\mathbf{w}} p_{\text{model}}(\mathbf{X}; \mathbf{w}) = \arg \max_{\mathbf{w}} \prod_{i=1}^N p_{\text{model}}(\mathbf{x}_i; \mathbf{w}) \\
 &= \arg \min_{\mathbf{w}} \sum_{i=1}^N -\log p_{\text{model}}(\mathbf{x}_i; \mathbf{w}) \\
 &= \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [-\log p_{\text{model}}(\mathbf{x}; \mathbf{w})] \\
 &= \arg \min_{\mathbf{w}} H(\hat{p}_{\text{data}}(\mathbf{x}), p_{\text{model}}(\mathbf{x}; \mathbf{w})) \\
 &= \arg \min_{\mathbf{w}} D_{\text{KL}}(\hat{p}_{\text{data}}(\mathbf{x}) \| p_{\text{model}}(\mathbf{x}; \mathbf{w})) + H(\hat{p}_{\text{data}}(\mathbf{x}))
 \end{aligned}$$

Maximum Likelihood Estimation

MLE can be easily generalized to the conditional case, where our goal is to predict t given \mathbf{x} :

$$\begin{aligned}
 \mathbf{w}_{\text{MLE}} &= \arg \max_{\mathbf{w}} p_{\text{model}}(\mathbf{t} | \mathbf{X}; \mathbf{w}) = \arg \max_{\mathbf{w}} \prod_{i=1}^N p_{\text{model}}(t_i | \mathbf{x}_i; \mathbf{w}) \\
 &= \arg \min_{\mathbf{w}} \sum_{i=1}^N -\log p_{\text{model}}(t_i | \mathbf{x}_i; \mathbf{w}) \\
 &= \arg \min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, t) \sim \hat{p}_{\text{data}}} [-\log p_{\text{model}}(t | \mathbf{x}; \mathbf{w})] \\
 &= \arg \min_{\mathbf{w}} H(\hat{p}_{\text{data}}(t | \mathbf{x}), p_{\text{model}}(t | \mathbf{x}; \mathbf{w})) \\
 &= \arg \min_{\mathbf{w}} D_{\text{KL}}(\hat{p}_{\text{data}}(t | \mathbf{x}) || p_{\text{model}}(t | \mathbf{x}; \mathbf{w})) + H(\hat{p}_{\text{data}}(t | \mathbf{x}))
 \end{aligned}$$

where the conditional entropy is defined as $H(\hat{p}_{\text{data}}) = \mathbb{E}_{(\mathbf{x}, t) \sim \hat{p}_{\text{data}}} [-\log(\hat{p}_{\text{data}}(t | \mathbf{x}))]$ and the conditional cross-entropy as $H(\hat{p}_{\text{data}}, p_{\text{model}}) = \mathbb{E}_{(\mathbf{x}, t) \sim \hat{p}_{\text{data}}} [-\log(p_{\text{model}}(t | \mathbf{x}; \mathbf{w}))]$.

The resulting *loss function* is called **negative log-likelihood (NLL)**, or **cross-entropy**, or **Kullback-Leibler divergence**.

An extension of perceptron, which models the conditional probabilities of $p(C_0|\mathbf{x})$ and of $p(C_1|\mathbf{x})$. Logistic regression can in fact handle also more than two classes, which we will see in the next lecture.

Logistic regression employs the following parametrization of the conditional class probabilities:

$$\begin{aligned}p(C_1|\mathbf{x}) &= \sigma(\mathbf{x}^T \mathbf{w} + b) \\p(C_0|\mathbf{x}) &= 1 - p(C_1|\mathbf{x}),\end{aligned}$$

where σ is a **sigmoid function**

$$\sigma(x) = \frac{1}{1 + e^{-x}}.$$

It can be trained using the SGD algorithm.

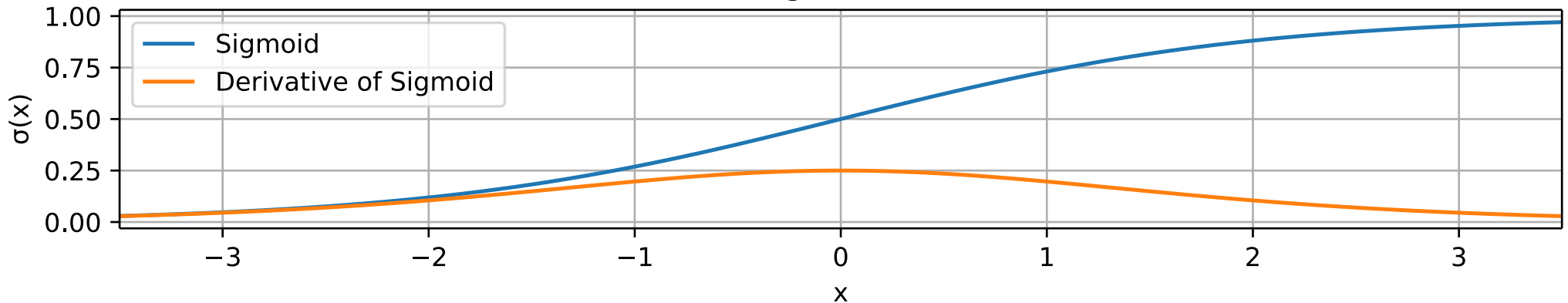
Sigmoid Function

The sigmoid function has values in range $(0, 1)$, is monotonically increasing and it has a derivative of $\frac{1}{4}$ at $x = 0$.

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Plot of the Sigmoid Function $\sigma(x)$



We denote the output of the “linear part” of the logistic regression as

$$\bar{y}(\mathbf{x}; \mathbf{w}) = \mathbf{x}^T \mathbf{w},$$

and the overall prediction as

$$y(\mathbf{x}; \mathbf{w}) = \sigma(\bar{y}(\mathbf{x}; \mathbf{w})) = \sigma(\mathbf{x}^T \mathbf{w}).$$

Logistic Regression

To train the logistic regression, we use MLE (the maximum likelihood estimation). Its application is straightforward, given that $p(C_1|\mathbf{x}; \mathbf{w})$ is directly the model output $y(\mathbf{x}; \mathbf{w})$.

Therefore, the loss for a minibatch $\mathbb{X} = \{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N)\}$ is

$$E(\mathbf{w}) = \frac{1}{N} \sum_i -\log(p(C_{t_i}|\mathbf{x}_i; \mathbf{w})).$$

Input: Input dataset $(\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{t} \in \{0, +1\}^N)$, learning rate $\alpha \in \mathbb{R}^+$.

- $\mathbf{w} \leftarrow \mathbf{0}$ or we initialize \mathbf{w} randomly
- until convergence (or patience runs out), process a minibatch of examples \mathbb{B} :
 - $\mathbf{g} \leftarrow \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \nabla_{\mathbf{w}} \left(-\log(p(C_{t_i}|\mathbf{x}_i; \mathbf{w})) \right)$
 - $\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{g}$

Everything we learned about **features** and L^2 **regularization** holds for logistic regression too.



After this lecture you should be able to

- Think about binary classification using **geometric intuition** and use the **perceptron algorithm**.
- Define the **main concepts of information theory** (entropy, cross-entropy, KL-divergence) and prove their properties.
- Derive training objectives using the **maximum likelihood principle**.
- Implement and use **logistic regression** for binary classification with SGD.